

KSH Államigazgatási Számítógépes Szolgálat

Információelméleti módszeren alapuló lényegkiemelési eljárás és klinikai adatbázis redukciója

Balogh Gábor, Götl Győző és Srajber Benedek

Bevezetés

Az orvostudomány és a pszichológia gyakran használt vizsgálati eszközei a tesztek, amelyekkel a vizsgált tényezők vagy tényezőcsoportok egymásrahatásának mértékére, illetve fontosságuk sorrendjére vonhatunk le következtetéseket. Mivel azonban a tesztek kiértékelése során fellépő valószínűségi változók értékkészletei (azaz a kérdésekre adható válaszok) rendezetlen, diszkrét halmazok, ezért a statisztika klasszikus módszerei (korrelációszámítás, regresszióanalízis stb.) csak igen korlátozottan alkalmazhatók. A rendezetlen halmazokat rendezettekre leképezve ugyanis korrelációs vizsgálatokkal csak a lineáris függés erősségét mérhetjük. (Az eredmény természetesen a leképezéstől is függ.)

E problémák megoldására születtek meg a statisztika újabb fejezetei és az alakfelismerés statisztikus módszerei: lényegkiemelés, hierarchikus- és szekvenciális Cluster-analízis stb.

Az általunk használt információs mérték diszkrét valószínűségi változók egymástól való függése tényleges nagyságának meghatározására szolgál. Mivel ennek kiszámításához csak a változók együttes eloszlását leíró ún. kontingencia-táblázat szükséges, emiatt attól is független, hogy a változókat nominális-, ordinális-, intervallum- vagy arányskálával jellemezzük. Az információs mérték képezi az alapját adatbázis-redukciós algoritmusunknak is.

1. Információs mérték

Legyen Z tetszőleges diszkrét, véges értékkészletű

$$p_i = p(Z = z_i) \quad /i = 1, \dots, n/ \text{ eloszlásu}$$

$$\sum_{i=1}^n p_i = 1$$

valószínűségi változó. Ekkor Z kimenetelének bizonytalanságát entrópiája, azaz a

$$H(Z) = \sum_{i=1}^n p_i \log \frac{1}{p_i} \quad \text{mennyiség méri.}$$

Ha Z vektorváltozó, azaz $Z = (X, Y)$, ahol X és Y valószínűségi változók, akkor (a Jensen-egyenlőtlenség alkalmazásával) a

$$H(Z) = H(X, Y) \leq H(X) + H(Y)$$

egyenlőtlenséghez jutunk. Az egyenlőség csak X és Y függetlensége esetén teljesül /lásd (3)/.

Az egyenlőtlenség két oldala közti különbséget kölcsönös információnak nevezzük és $I(X, Y)$ -nal jelöljük. Az elnevezést az indokolja, hogy az X változó éppen $I(X, Y)$ bit információt tartalmaz az Y -ra vonatkozóan (és ugyanennyit tartalmaz Y az X -re nézve).

Tehát az $I(X, Y) = H(X) + H(Y) - H(X, Y)$ bevezetésével az $I(X, Y) \geq 0$. Ismeretes, hogy $I(X, Y) \leq H(Y)$ ill. $I(X, Y) \leq H(X)$ is fennáll, továbbá pl. az $I(X, Y) \leq H(X)$ -ben pontosan akkor van egyenlőség, ha X az Y függvénye.

Következésképpen az $i(X, Y) = \frac{I(X, Y)}{H(X)}$ formulával meghatározott mennyiségre $0 \leq i(X, Y) \leq 1$, $i(X, Y) = 1$ akkor és csak akkor, ha X függvénye Y -nak. Így $i(X, Y)$ olyan nem-szimmetrikus mennyiség, mely X -nek Y -tól való függését méri, származtatása alapján információs mértéknek nevezzük.

Természetesen, két valószínűségi változó kapcsolatának részletesebb vizsgálatánál nem lehet elegendő csupán az információs mérték és a korrelációs értékének meghatározása. Ezek ugyanis csak a két változó globális kapcsolatát jellemzik. Azt, hogy az X változó egy konkrét kimenetele milyen mértékben határozza meg Y -t, az $i(Y, X=x_i)$ parciális információs mérték mutatja meg:

$$i(Y, X=x_i) = \frac{H(Y) - H(Y, X=x_i)}{H(Y)}$$

Könnyen belátható, hogy

$$i(Y, X) = \sum_{j=1}^n p(X=x_j) \cdot i(Y, X=x_j) \text{ és } 1 - \frac{\log k}{H(Y)} \leq i(Y, X=x_j) \leq 1.$$

Ha $i(Y, X=x_j) = 1$, akkor az $X=x_j$ kimenetel teljes mértékben meghatározza Y kimenetelét, ha $i(Y, X=x_j) = 0$, akkor bizonytalanságunk változatlan marad, ha pedig negatív, akkor a szóban forgó megfigyelés növeli bizonytalanságunkat.

Legyenek most X_i -k $/i = 1, \dots, k/$ és Y valószínűségi változók.

Igazolható, hogy

$$\begin{aligned} I(Y, X_1, \dots, X_k) - I(Y, X_1, \dots, X_{k-1}) &= \\ &= I(X_k, X_1, \dots, X_{k-1}, Y) - I(X_k, X_1, \dots, X_{k-1}), \end{aligned} \quad /1/$$

amely azt mutatja, hogy X_k kihagyásával Y -ról annyi információt veszünk, amennyit X_k -ről nyerhetünk Y megismerése által. Az előző egyenlőség segítségével kapható az

$$i(Y, X_1, X_2) \geq i(Y, X_1) + i(Y, X_2) - \frac{I(X_1, X_2)}{H(Y)}$$

egyenlőtlenség. Egyenlőség akkor és csak akkor van, ha

$$i(X_2, Y, X_1) = i(X_2, Y),$$

ami azzal ekvivalens, hogy " X_1 és X_2 független Y bármely kimenete esetén". (Megjegyezzük, hogy ennek semmi kapcsolata a valószínűségi változók között értelmezett függetlenség fogalmával, amint azt Lee (2) munkájában feltételezte.) Így független valószínűségi változók esetén is csak az $i(Y, X_1, X_2) \geq i(Y, X_1) + i(Y, X_2)$ egyenlőtlenség igaz.

Az információs mérték közelítésével és a közelítés szignifikanciájának vizsgálatával az (1) tanulmány foglalkozik. Ezeket a problémákat a következőképpen foglalhatjuk össze:

Jelöljük két valószínűségi változó információs mértékének elméleti értékét $i(Y, X)$ -szel, gyakorlati meghatározására N számú kísérletet végzünk. Azaz N elemű mintapopuláció alapján meghatározzuk az (Y, X) tapasztalati eloszlásfüggvényét és ennek segítségével kiszámítjuk az ehhez a mintakollekcióhoz tartozó tapasztalati információs mértéket,

$$i_{\text{tap}}^N(Y, X) \text{ -et.}$$

Belátható, hogy a mintaanyag nagyságának növekedésével a tapasztalati információs mérték sztochasztikusan konvergál az információs mérték elméleti értékéhez.

II. Adatredukciós algoritmusunk elméleti megfontolása

Az előző fejezetben megtárgyaltuk az információs mérték tulajdonságait és utaltunk becslésének módjára is. Ebben a részben rámutatunk az információs mérték fontosságára a lényegkiemelő és adatredukciós módszerek egy osztályában.

Az Y és az X , véges számú diszkrét értéket felvevő valószínűségi változók függetlenségét vizsgáljuk.

Null-hipotézisünk H_0 : Y és X függetlenek,

ellen-hipotézisünk H_e : nem-függetlenek

$$\{p(Y_i) \cdot p(x_i)\} \quad \text{illetve} \quad \{p(Y_i, x_i)\} \quad \text{eloszlással.}$$

Döntés céljából n számú független megfigyelést végzünk. Egy megfigyeléssorozat eredményét az

$$u = (Y_{i_1}, x_{i_1}, Y_{i_2}, x_{i_2}, \dots, Y_{i_n}, x_{i_n})$$

jellemzi. A megfigyelések függetlensége miatt természetesen:

$$p^n(u | H_0) = \prod_{k=1}^n p(Y_{i_k}) \cdot p(x_{i_k}) \quad \text{ill.} \quad p^n(u | H_e) = \prod_{k=1}^n p(Y_{i_k}, x_{i_k}).$$

Jelöljük az n megfigyelés összes lehetséges kimeneteleinek halmazát R -rel és legyen E (elfogadási tartomány) R -nek egy részhalmaza.

Induljunk ki abból az elvből, hogy $u \in E$ esetén H_0 , $u \notin E$ esetén pedig H_e javára döntünk. Célunk akkor az E halmaz olyan meghatározása, hogy a helytelen döntésből eredő hiba minél kisebb legyen. Tehát olyan E halmazt keresünk, melynél az elsőfajú

$$p^n(E^c | H_0) = \sum_{u | u \notin E} p^n(u | H_0),$$

ill. a másodfajú

$$p^n(E | H_e) = \sum_{u | u \in E} p^n(u | H_e)$$

hibák minimálisak. Mivel ezek változása ellentétes irányultságú, azért minimumuk egyszerre nem valósítható meg. Ehelyett - rögzített esetén - vizsgálhatjuk a következő kifejezés n -től való függését:

$$\chi(n) = \inf p^n(E | H_e),$$

$E: ECR$ és

$$p^n(E | H_0) \geq 1 - \beta$$

azaz azt, hogy az elsőfajú hiba valószínűségét szignifikancia-szinten tartva hogyan függ a minimális másodfajú hiba valószínűsége a mintaadanyag nagyságától. Erre vonatkozik a Chernofftól származó tétel (6):

$$\lim_{n \rightarrow \infty} \frac{1}{n} \cdot \log(n) = -I(Y, X).$$

Tehát, az elsőfajú hiba valószínűségét szignifikancia-szinten tartva a másodfajú hiba valószínűsége a mintaszám növelésével exponenciálisan csökkenthető. A csökkenés mértékét az $I(Y, X)$ kölcsönös információ határozza meg. Így, ha X vektorváltozó, azaz $X = (X_1, \dots, X_s)$, akkor adatredukálás céljából azon X_{r_1}, \dots, X_{r_m} változók kiválasztása kívánatos, melyekre az $I(Y, X_{r_1}, \dots, X_{r_m})$ kölcsönös információ és ebből következően az $i(Y, X_{r_1}, \dots, X_{r_m})$ információs mérték maximális.

III. Adatredukciós algoritmus

A feladat azon minimális számú X_{r_1}, \dots, X_{r_m} változóknak az X_1, \dots, X_s változók közül történő kiválasztása, amelyek Y -ról az összes változó által nyerhető információnak "elég nagy részét" tartalmazzák. Pontosabban: keressük azokat az X_{r_1}, \dots, X_{r_m} változókat, amelyekre

$$\left\{ X_{r_j} \right\}_{j=1}^m \subset \left\{ X_i \right\}_{i=1}^s \quad \text{és}$$

$i(Y, X_{r_1}, \dots, X_{r_m}) \geq (1 - \varepsilon) \cdot i(Y, X_1, \dots, X_s)$ feltételek mellett m értéke minimális.

A feladat egzakt megoldhatóságának két fő akadálya van:

- 1.) $i(Y, X_1, \dots, X_s)$ gyakorlatilag kiszámíthatatlan, értékét csak becsülni lehet.
- 2.) Az összes lehetséges változó-kombinációk leszámhlálása gyakorlatilag elvégezhetetlen.

Ezen okok miatt nem is várható, hogy létezik optimális megoldáshoz vezető, gyors algoritmus. A következőkben leírt "ad hoc" módszer célja ezért nem az optimális, hanem egy ahhoz nagy valószínűséggel közel eső megoldás keresése.

Iterációs algoritmusunk lépései:

- 1.) Kiválasztjuk az $i(Y, X_{r_1}) = \max_{1 \leq j \leq s} i(Y, X_j)$ feltételnek eleget tevő X_{r_1} -et.
- 2.) Meghatározzuk az $i(Y, X_{r_1}, X_{r_2}) = \max_{\substack{1 \leq j \leq s \\ j \neq r_1}} i(Y, X_{r_1}, X_j)$ feltételt kielégítő X_{r_2} -t.
- ...
- k.) Az előző lépésekben meghatározott $X_{r_1}, \dots, X_{r_{k-1}}$ változókat rögzítve

meghatározzuk azt az újabb X_{r_k} változót, amelyre

$$i(Y, X_{r_1}, \dots, X_{r_k}) = \max_{\substack{1 \leq j \leq s \\ i \neq r_q / q = 1, \dots, k-1}} i(Y, X_{r_1}, \dots, X_{r_{k-1}}, X_j) \quad /2/$$

teljesül. Természetesen az információs mérték értékét az 1. részben leírtak alapján tapasztalati értékével közelítjük.

A befejezhetőség kritériumai a következők:

- a.) A kiválasztott változók által Y -ról nyert információ eléggé megközelíti-e az összes kapható információt?
- b.) Feltételezhető-e, hogy újabb változó kiválasztásával még további szignifikáns mennyiségű információt nyerünk?

Az algoritmus számítógépes realizálásával kapcsolatban szükségesnek tartjuk megjegyezni a következőket:

a.) A lépésenkénti számolás mennyisége a lépésszám növekedésével exponenciálisan nő, ezért a kiválasztható változók száma elsősorban a mintaanyag nagyságától és a feldolgozó számítógép paramétereitől függ.

b.) Az /1/ összefüggés figyelembevétele alapján az algoritmus-sal kiválasztott változók egymástól csak kevésbé függhetnek, mivel módszerünk csak akkor választ ki egy, az előzőktől erősen függő változót, ha másképpen nem kaphatnánk Y -ról legalább ugyanannyi információt.

c.) A megoldhatóság nehézségei miatt a befejezést eldöntő vizsgálatokat függvény approximációs módszerekkel végezhetjük el.

d.) Az algoritmus részletes leírása és FORTRAN programja az ÁSZSZ Honeywell 66/60-as gépének programkönyvtárában az érdeklődők rendelkezésére áll. A program próbafeladaton sikeresen lefutott. Nagyméretű mintaanyagon kipróbálása folyamatban van.

IV. Klinikai adatbázis redukciója, differenciál-diagnózis és prognózis

A diagnosztikai folyamat végrehajtása közben az orvos rendszert szembetalálja magát a nehézséggel, hogy a nagy számosságú evidencia halmaz nehezen kezelhető vagy egyszerűen áttekinthetetlen számára. Ebből a tényből hármasszámú probléma adódik:

a.) Az evidencia halmaz elemeinek feltétlen szükséges redukciója

Ezt megoldhatjuk oly módon, hogy a III. pontban ismertetett algoritmus /2/ képletébe Y -nak az $Y = (X_1, \dots, X_s)$ valószínűségi vektorváltozót választjuk, ahol X_i -k ($i=1, \dots, s$) mint valószínűségi változók, a beteg adatainak felelnek meg. Az eredmény: a betegre vonatkozó tényezők számának csökkenése oly módon, hogy a redukált tényezők halmaza közel azonos információ-mennyiséget tartalmaz, mint a kiindulásul vett tényezők együttesen.

b.) A differenciál-diagnózis megállapításához szükséges legszűkebb evidencia halmaz megállapítása

Ha egy klinikai adatbázis esetén a /2/ képletben Y a szóbaeső betegségeket osztályokat reprezentálja, akkor az ismertetett adatredukciós eljárás az összes vizsgált tényező közül azon minimális számú "közel független" tényező kiválasztását eredményezi, amely elegendő információt tartalmaz a betegség-osztályok valószínűsítéséhez. Kevés számú tényező esetén most már a betegek betegség-osztályba sorolásához a jól ismert Bayes-féle differenciál-diagnosztikai modellel (9) vagy más módszerekkel is könnyebben eljuthatunk.

c.) A beteg prognózisához (pl. túlélések valószínűsítése) szükséges legszűkebb halmaz meghatározása

Ha a /2/ képletben Y - mint valószínűségi változó - valamely betegség összes lehetséges kimeneteleit, X_i -k $/i=1, \dots, s/$ pedig a kimenetelt befolyásoló összes szóbaeső tényezőt reprezentálják, az eredmény b.)-hez (adatredukció és osztályba sorolás) hasonlóan a kívánt redukált halmaz és a megfelelő kimeneteli kategóriába való besorolás.

Irodalom

- (1) Andor Cs., Mérő L.: Kölcsönös információ alkalmazása két vagy több változó kapcsolatának meghatározására. MRT TCI közl.
- (2) R.Char-Tung Lee: Application of Information Theory to Relevant Variables. Math.Biosci. 11,153-161 (1971)
- (3) Csiszár I., Fritz J.: Információelmélet, 1970.
- (4) Fritz J.: Bevezetés az információelméletbe (Tankönyvkiadó)
- (5) J.Kryspin, A.M.Norwich: Use of Information Theory in Analysis of Medical Data. 1972 IEEE Conference on Inf. Theory, Asilomar, California.
- (6) S.Kullback: Information Theory and Statistics. Wiley, 1959.
- (7) W.S.Meisel: Computer-oriented Approaches to Pattern Recognition. Acad. Press, New York, 1972.
- (8) J.M. Mendel, K.S.Fu: (szerk) Adaptive, Learning and Pattern Recognition System. Acad. Press, 1970.
- (9) L.G. Withby, W.Lutz: Principles and Practice of Medical Computing.

